

# *Gene finding and Genome annotation*

Manfred Zorn

**BerkeleyPGA**

Bioinformatics Tools for Comparative Analysis

July 16, 2002

## What is a Gene?

- **Definition:** An inheritable trait associated with a region of DNA that codes for a polypeptide chain or specifies an RNA molecule which in turn have an influence on some characteristic phenotype of the organism.

Abstract concept that describes a complex phenomenon

## What is Annotation?

- **Definition:** Extraction, definition, and interpretation of features on the genome sequence derived by integrating computational tools and biological knowledge.

Identifiable features in the sequence

## How does an annotation differ from a gene?

- Many annotations describe features that constitute a gene.
- Other annotations may not always directly correspond in this way, e.g., an STS, or sequence overlap



## DNA Analysis

---

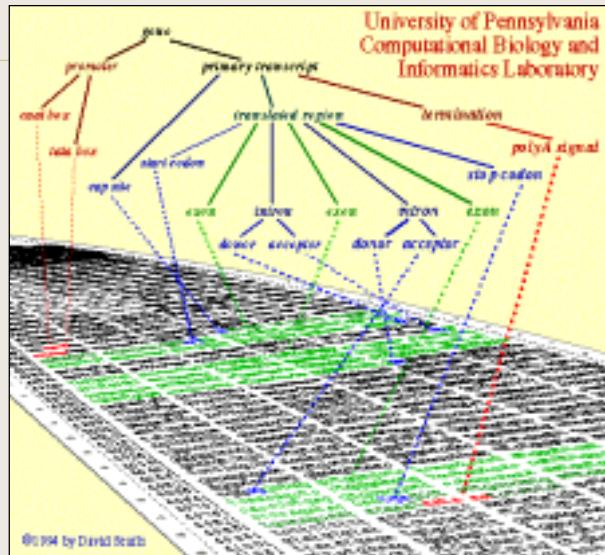
- Heuristics
- Statistics
- Artistics

## DNA Analysis

---

- Find the genes
  - Heuristic signals
  - Inherent features
  - Intelligent methods
- Characterize each gene
  - Compare with other genes
  - Find functional components
  - Predict features

# What is a Gene?

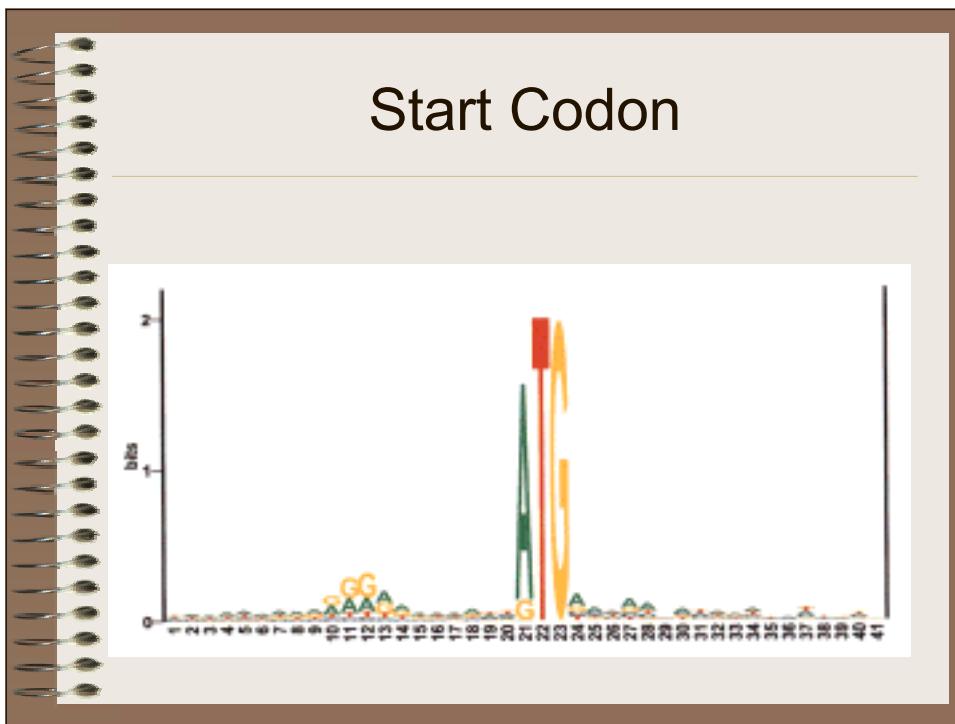


## Heuristic Signals

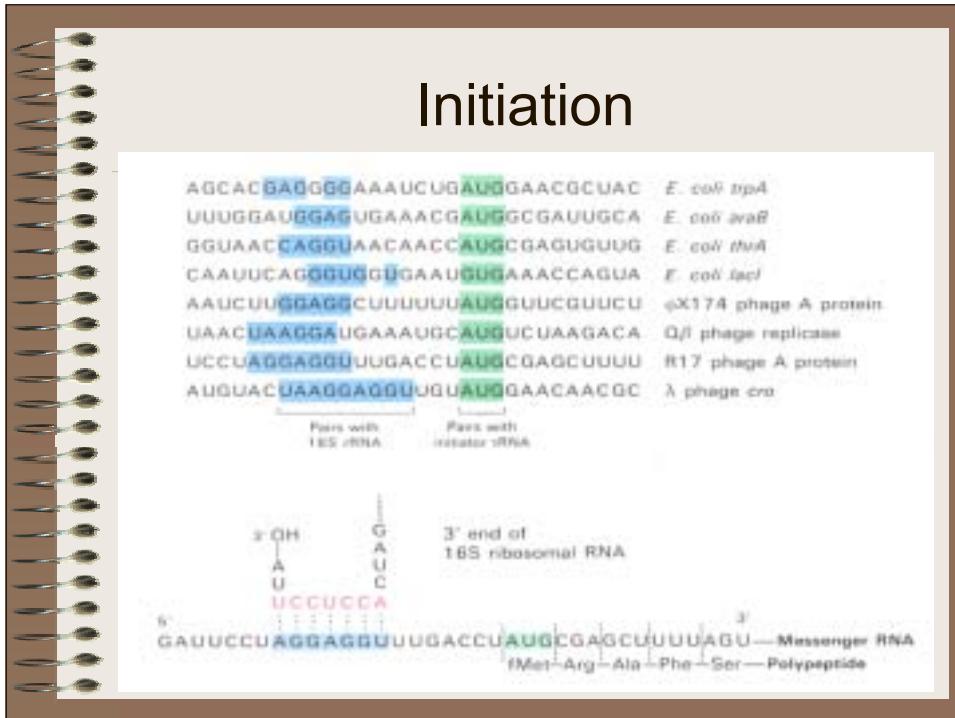
- DNA contains various recognition sites for internal machinery
- Promoter signals
- Transcription start signals
- Start Codon
- Exon, Intron boundaries
- Transcription termination signals



## Start Codon



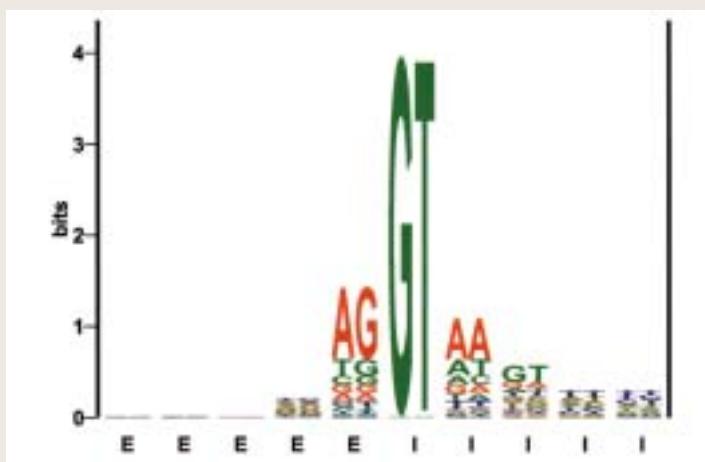
## Initiation



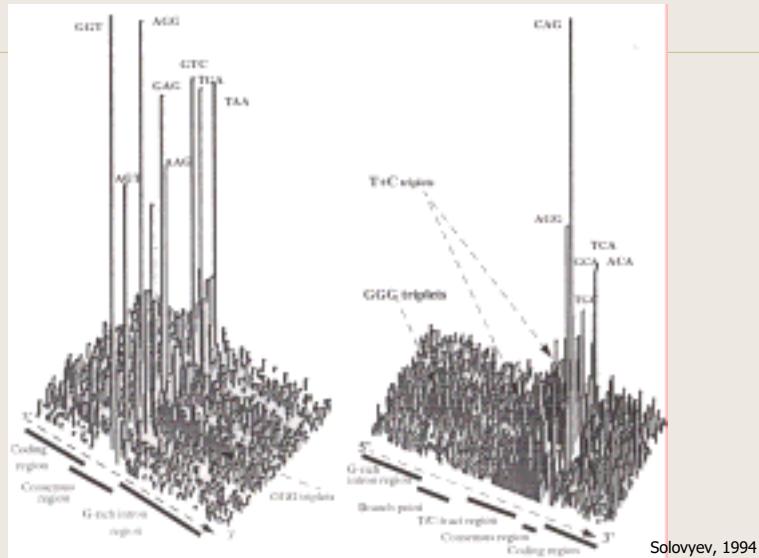
## Inherent Features

- DNA exhibits certain biases that can be exploited to locate coding regions
- Uneven distribution of bases
- Codon bias
- CpG islands
- In-phase words
- Encoded amino acid sequence
- Imperfect periodicity
- Other global patterns

## Donor Splice Site



## Inherent Features



Solovyev, 1994

## Intelligent Methods

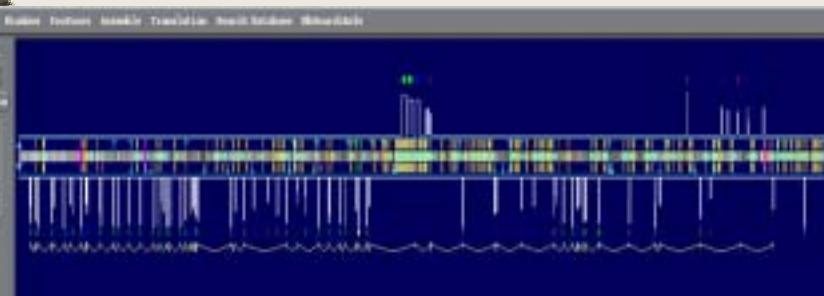
- Pattern recognition methods weigh inputs and predict gene location
  - Content-based methods
  - Site-based methods
  - Comparative methods
- Neural Networks
- Hidden Markov Models
- Stochastic Context-Free Grammar

## GRAIL *Uberbacher, Mural*

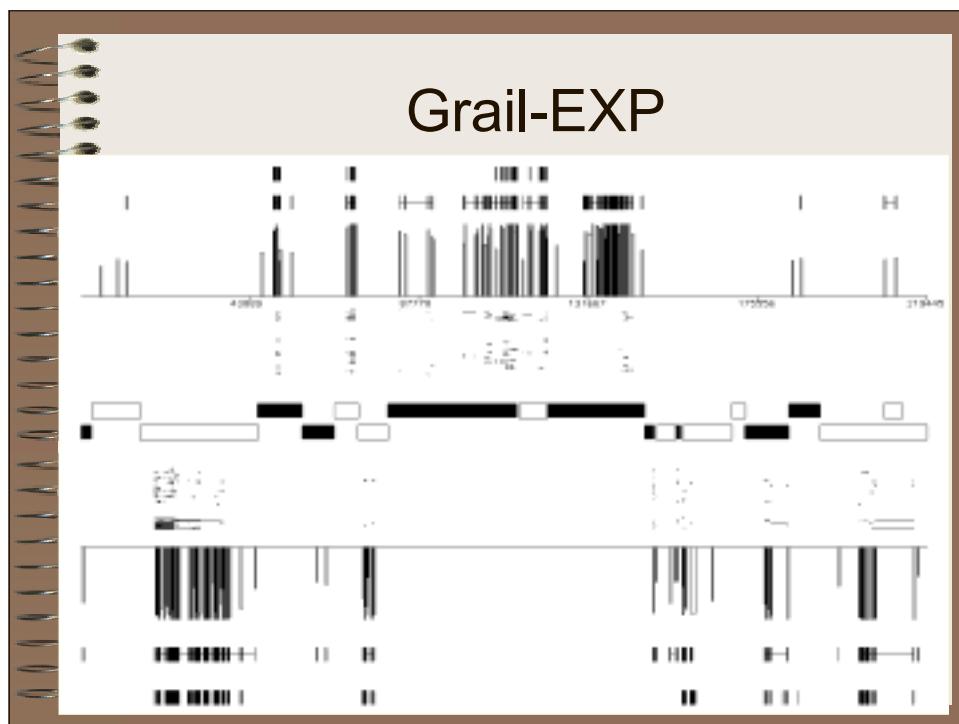
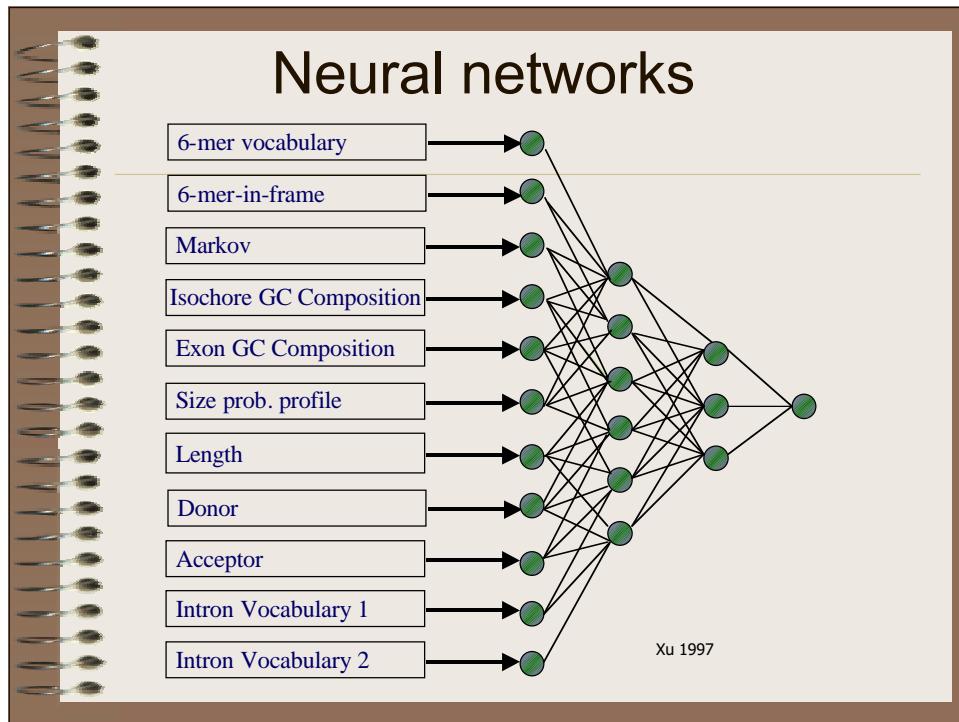
- GRAIL 1
  - Neural network with fixed window length (100 bases)
- GRAIL 1a
  - GRAIL 1 + adjacent information
- GRAIL 2
  - Variable length window, contextual information
- GRAIL-EXP
  - Comparison with partial and complete gene sequences

## Analyzing Complex Multi-Gene Regions

- Errors in exon prediction and splice site boundaries
- Gene boundaries uncertain
- Genes can be on both strands



Uberbacher



## FGNEH/FGENES Solovyev

- Looks at several structural features
  - Splice donor/acceptor sites
  - Putative coding regions
  - Intronic regions
- *Linear discriminant analysis* to split exon / non-exon classes
- Dynamic programming to assemble best gene structure

## MZEF Zhang

- *Quadratic discriminant analysis*
  - Exon length
  - Exon-intron transitions
  - Splice sites
  - Branch sites
  - Exon, strand, frame scores
- Detects internal exons
- No information about gene structure



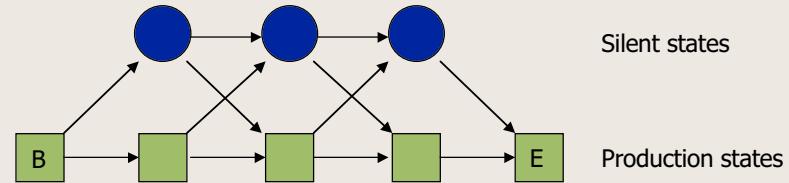
## GENSCAN *Burge, Karlin*

- Probabilistic model of sequence composition and gene structure
  - Looks for gene structure descriptions that are consistent with the query sequence to assign probability that sequence stretch is exon, ...
  - Best ---> optimal
  - But generates also suboptimal exons

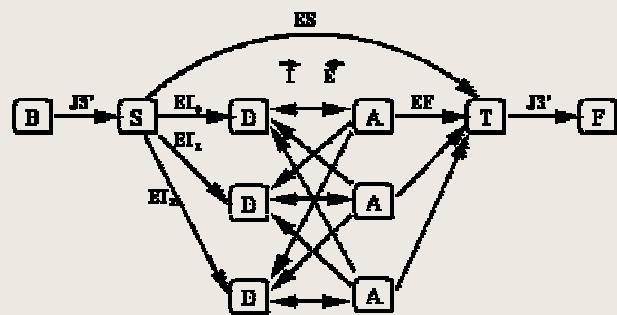
## PROCRUSTES *Gelfand*

- Forces sequence into target structure
  - Requires putative gene product
  - Stretches/shortens sequence to fit into model

## Hidden Markov Models



## GENIE *Kulp, Reese, Haussler*



## Strategies

- Select by correlation coefficient
- Select by review paper
- Select by recommendation
- Use them all

## Drawbacks

- Most programs are “trained” on existing data
- It’s awfully hard to find new things this way!
  - NTT
  - IPW

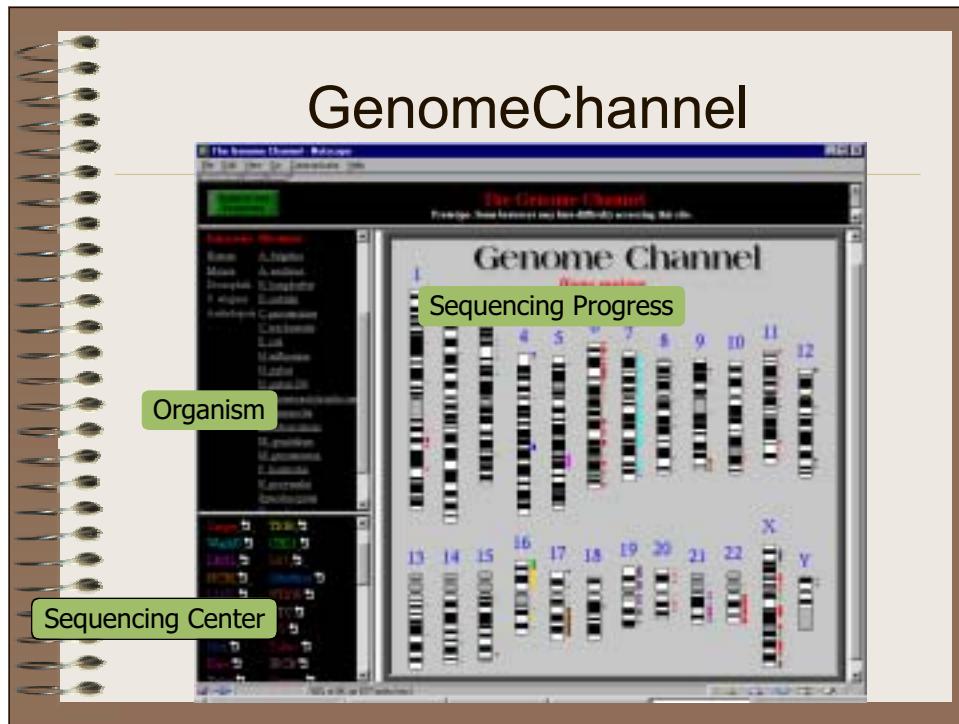
## Internet Resources

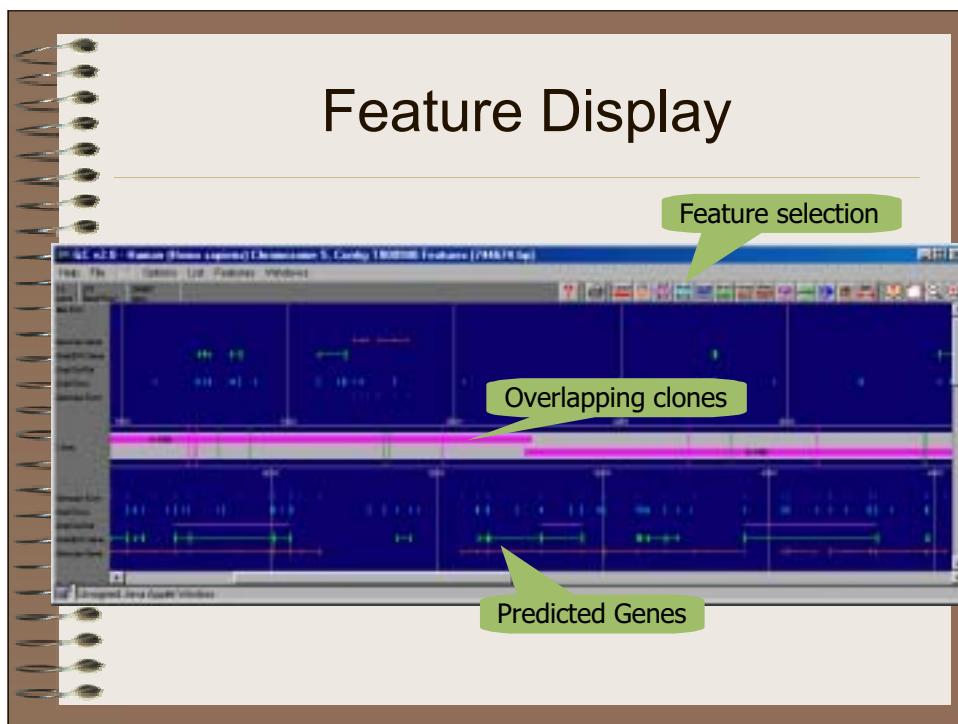
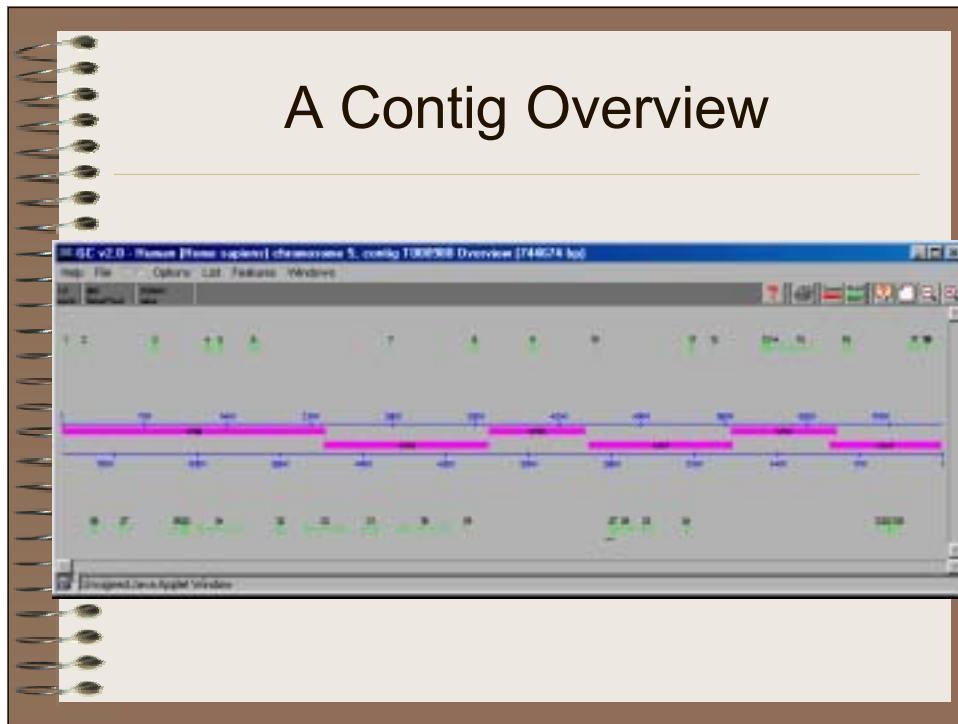
Banbury Cross	<a href="http://igs-server.cnrs-mrs.fr/igs/banbury">http://igs-server.cnrs-mrs.fr/igs/banbury</a>
FGENEH	<a href="http://genomic.sanger.ac.uk/gf/gf.shtml">http://genomic.sanger.ac.uk/gf/gf.shtml</a>
GeneID	<a href="http://www1.imim.es/geneid.html">http://www1.imim.es/geneid.html</a>
GeneMachine	<a href="http://genome.nhgri.nih.gov/genemachine">http://genome.nhgri.nih.gov/genemachine</a>
GENSCAN	<a href="http://genes.mit.edu/GENSCAN.html">http://genes.mit.edu/GENSCAN.html</a>
Genotator	<a href="http://www.fruitfly.org/_nomi/genotator/">http://www.fruitfly.org/_nomi/genotator/</a>
GRAIL	<a href="http://compbio.ornl.gov/tools/index.shtml">http://compbio.ornl.gov/tools/index.shtml</a>
GRAIL-EXP	<a href="http://compbio.ornl.gov/grailexp">http://compbio.ornl.gov/grailexp</a>
MZEF	<a href="http://www.cshl.org/genefinder">http://www.cshl.org/genefinder</a>
PROCRUSTES	<a href="http://www-ho.usc.edu/software/procrustes">http://www-ho.usc.edu/software/procrustes</a>
RepeatMasker	<a href="http://ftp.genome.washington.edu/RM/RepeatMasker.html">http://ftp.genome.washington.edu/RM/RepeatMasker.html</a>
HMMgene	<a href="http://www.cbs.dtu.dk/services/HMMgene">http://www.cbs.dtu.dk/services/HMMgene</a>
Chapter 10	<a href="http://www.wiley.com/legacy/products/subject/life/bioinformatics/chapterlinks.html">http://www.wiley.com/legacy/products/subject/life/bioinformatics/chapterlinks.html</a>

## Characterize a Gene

Collect clues for potential function

- Comparison with other known genes, proteins
- Predict secondary structure
- Fold classification
- Gene Expression
- Gene Regulatory Networks
- Phylogenetic comparisons
- Metabolic pathways





## Gene Summary Report

The screenshot shows a "Gene Summary Report" window with the following details:

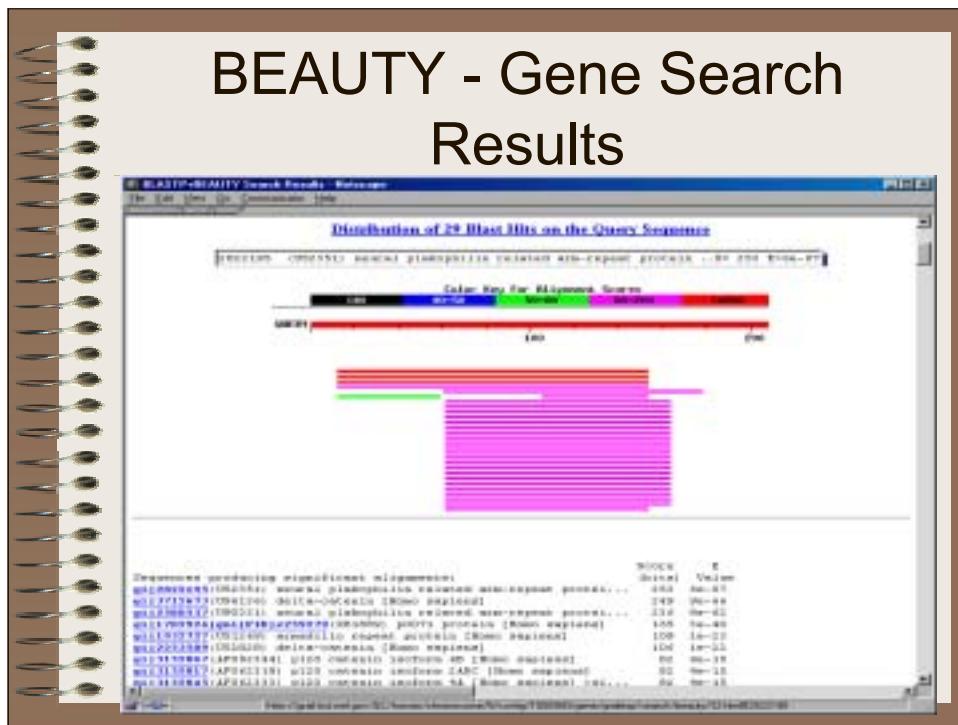
**Gene Information:**

- Type: gene
- ID: 04479\_pg
- Organism: Homo sapiens
- Chromosome: X
- Start: 87911
- End: 923341232
- Strand: -
- Protein ID: NP\_0012611
- Protein Name: C14orf111
- Protein Description: Zinc finger protein 111
- Exons: 12
- Transcripts: 1
- Gene ID: 123456789
- Gene Name: C14orf111
- Gene Description: Zinc finger protein 111
- Gene Ontology Terms: Biological Process: Gene expression, Molecular Function: Zinc finger domain binding, Cellular Component: Nucleoplasm
- Gene Structure: Shows exons 1-12 and their coordinates (e.g., 87911..87940, 923341232..923341261).
- Gene Expression: Shows mRNA levels (e.g., 0.000, 0.000, 0.000, 0.000, 0.000, 0.000, 0.000, 0.000, 0.000, 0.000, 0.000, 0.000).
- Gene Coverage: Shows coverage across the gene (e.g., 0.000, 0.000, 0.000, 0.000, 0.000, 0.000, 0.000, 0.000, 0.000, 0.000, 0.000, 0.000).

**Sequence Alignment:**

Sequence alignment details are present in the main panel, showing matches between the query sequence and the gene's coding region.

## BEAUTY - Gene Search Results

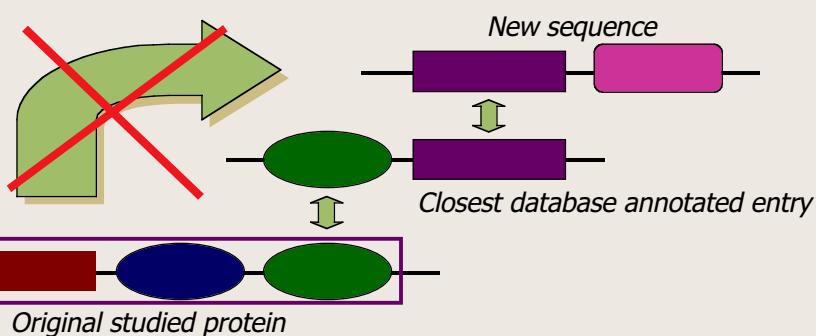


## Layers of Information

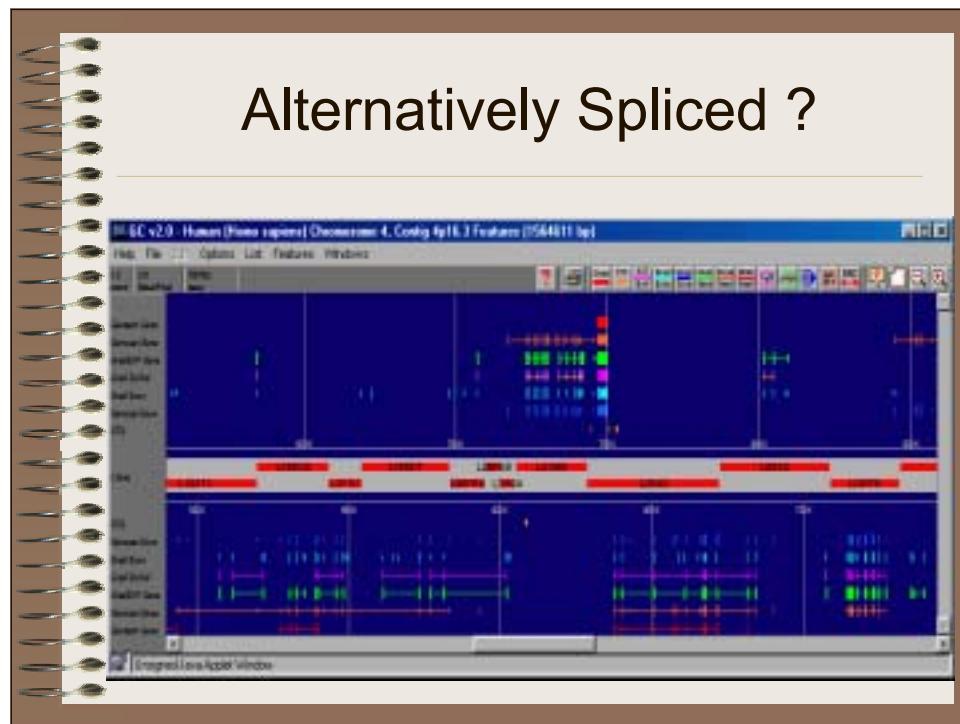
The same base sequence contains many layered instructions!

- Chromosome structure and function
  - Telomers, centromers
- Gene Regulatory information
  - Enhancers, promoters, ...
- Instructions for gene structure
- Instructions for protein
- Instructions for protein post-processing and localization

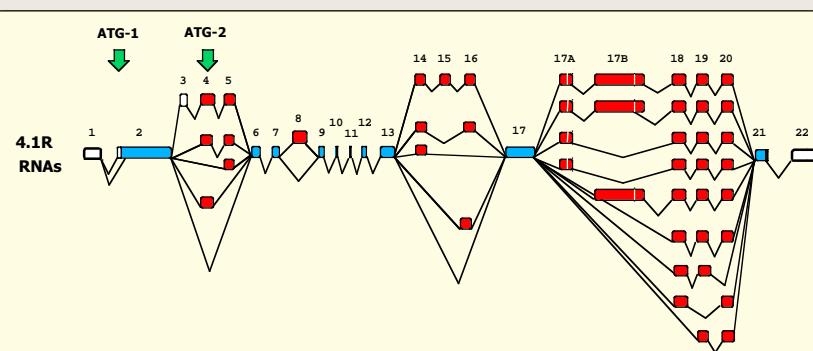
## Inherited Annotation Problems in Multi-Domain Proteins



## Alternatively Spliced ?

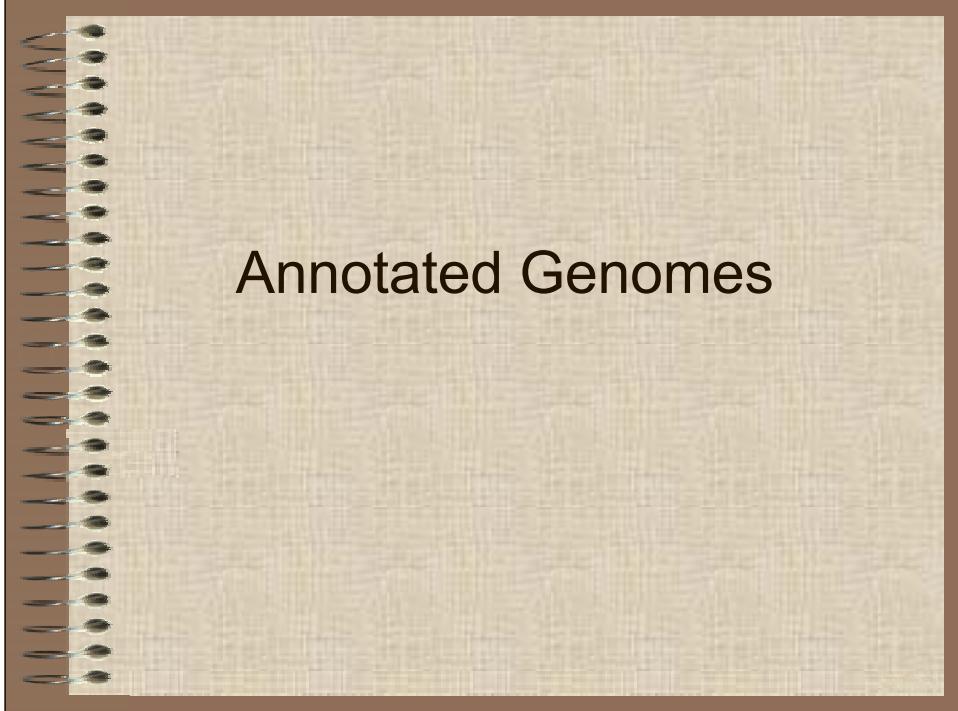


## One Gene - Many Proteins

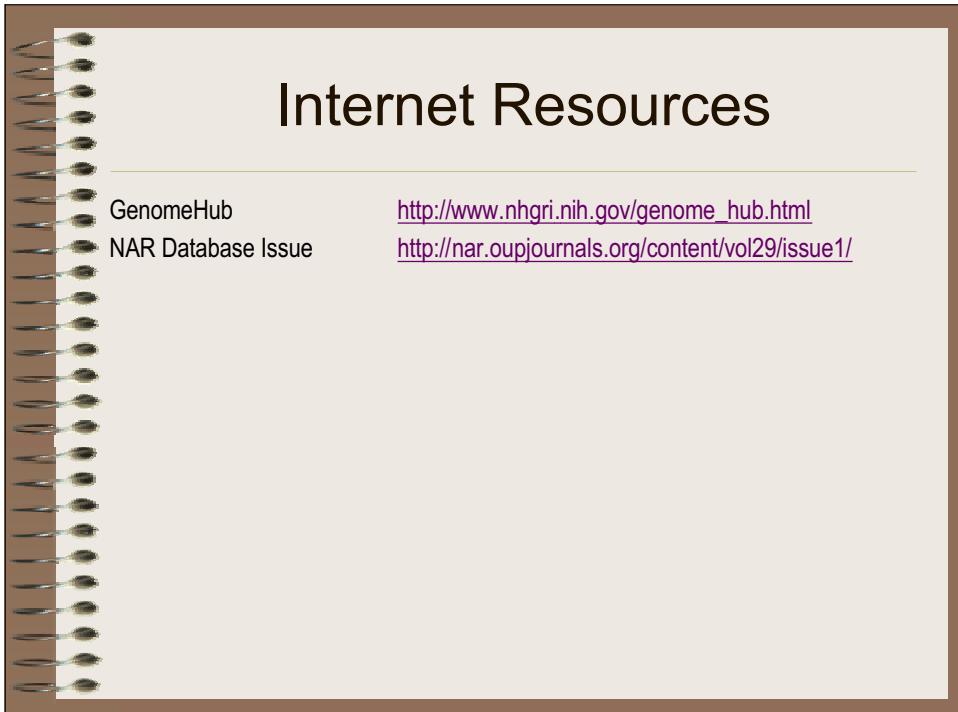


Conboy 1998

As many as 30% of human genes,  
in particular structural genes, may  
be alternatively spliced.



## Annotated Genomes



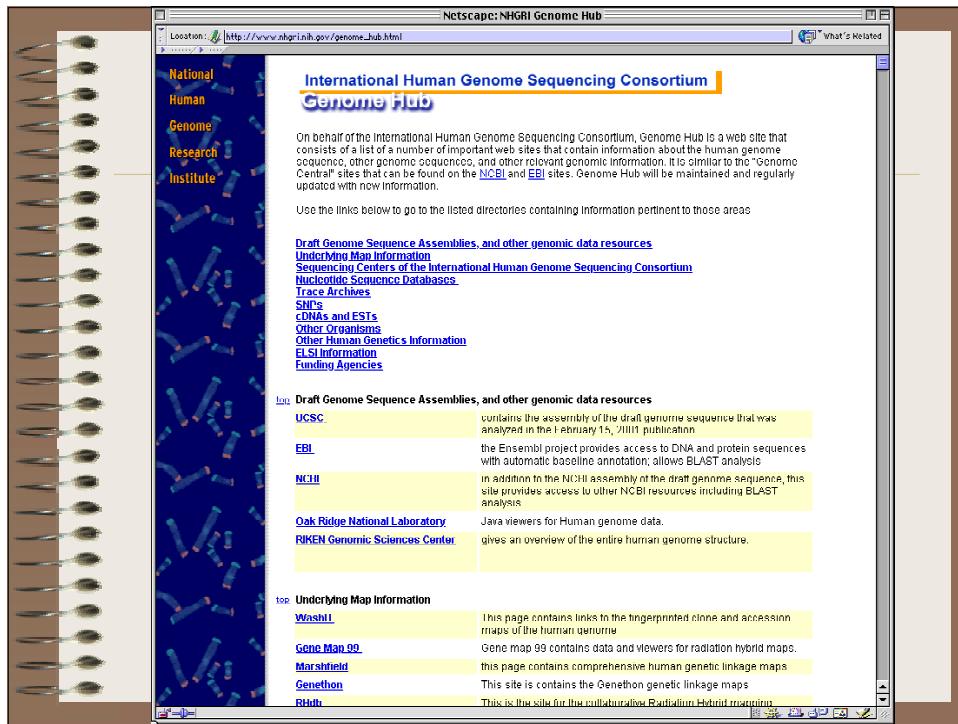
## Internet Resources

GenomeHub

[http://www.nhgri.nih.gov/genome\\_hub.html](http://www.nhgri.nih.gov/genome_hub.html)

NAR Database Issue

<http://nar.oupjournals.org/content/vol29/issue1/>



# GoldenPath

Netscape: Human Genome Browser Gateway

<http://genome.cse.ucsc.edu/goldenPath/octTracks.html>

**Human Genome Browser**

Web tool created by Jim Kent of UC Santa Cruz  
7 Oct 2000 draft assembly of the human genome

Please enter a position in the genome, set your preferred window width, and press the submit button. (Use BLAT Search to locate a particular sequence in the genome.)

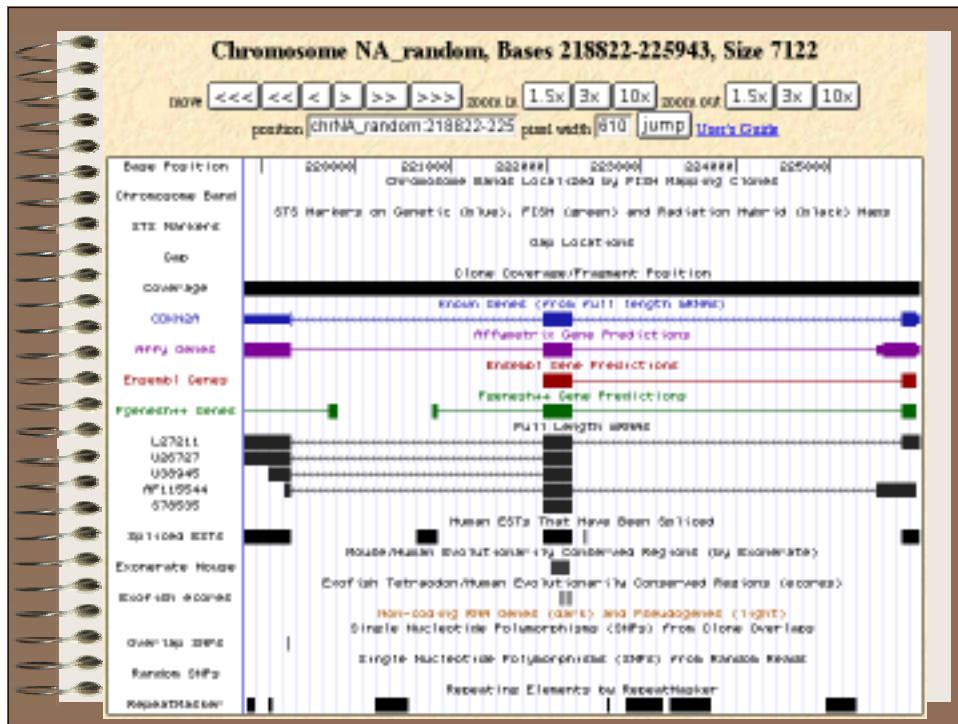
genome position:  pixel width:

A genome position can be specified by the accession number of a sequenced human genomic clone, an mRNA or EST or STS marker, or a cytological band, a chromosomal coordinate range, or keywords from the Genbank description of an mRNA. See the User Guide for more help.

**Request:** **Genome Browser Response:**

- chr19 Displays all of chromosome 19
- 20p13 Displays region for band p13 on chr 20
- 4q28 Displays band q28 on chr 4, gene region determining red hair color
- chr3-1:1000000 Displays first million bases of chr 3, counting from p arm telomere
- D16S3046 Displays region around STS marker D16S3046 from the Genethon/Marschall maps (open "STS Markers" track by clicking to see this marker)
- AA205474 Displays region of EST with GenBank acc. AA205474 in BRCA1 cancer gene on chr 17 (open "spliced ESTs" track by clicking to see this EST)
- ctg13698 Displays region of the fingerprint clone contig ctg13698 from the Wash. U. map (set "PVC contigs" track to "dense" and refresh to see PVC contigs)
- APU01670 Displays region of clone with GenBank accession APU01670 (open "coverage" track by clicking to see this clone)
- AF083811 Offers 2 coordinate choices on chr 7 and chr 22 for this mitotic checkpoint mRNA
- PRNP Offers coordinates choices for nested mRNAs for prion gene, chr 20p
- pectenope mRNA Lists transcribed pseudogenes but not cDNAs
- homeobox caudal Lists mRNAs for caudal homeobox genes
- vmyl tRNA Lists mRNAs for vmyl tRNA synthetase but not unsubmitted tRNA
- sinc finger Lists many sinc finger mRNAs
- kruppel zinc finger Lists only kruppel like sinc fingers
- huntington Lists candidate genes associated with Huntington's disease
- zohler Lists mRNAs deposited by scientist named Zohler
- Evans J.E. Lists mRNAs deposited by co author J.E. Evans

Use this list for entry authors -- even though Genbank searches require Evans J.E format, GenBank entries themselves use Evans, J.E. internally.



# BLAT Search

Netscape:BLAT Search

Location: <http://genome.cse.ucsc.edu/cgi-bin/hgDist?db=hg5> What's Related

BLAT Search Human Genome

From: Oct 7, 2000      Query type: BLAT's guess      Sort output: query\_score      Submit

Please paste in a query sequence to see where it is located in the UCSC assembly of the human genome. Multiple sequences can be searched at once if separated by a line starting with > and the sequence name.

Only DNA sequences less than 20,000 bases and protein or translated sequence of less than 4000 letters will be processed. If multiple sequences are submitted at the same time, the total limit is 50,000 bases or 10,000 letters.

BLAT on DNA is designed to quickly find sequences of 95% and greater similarity of length 40 bases or more. It may miss more divergent or shorter sequence alignments. It will find perfect sequence matches of 33 bases, and sometimes find them down to 22 bases. BLAT on proteins finds sequences of 80% and greater similarity of length 20 amino acids or more. In practice DNA BLAT works well on primates, and protein BLAT on land vertebrates.

BLAT is not BLAST. DNA BLAT works by keeping an index of the entire genome in memory. The index consists of all non-overlapping 11-mers except for those heavily involved in repeats. The index takes up a bit less than 2 gigabytes of RAM. The genome itself is not kept in memory, allowing BLAT to deliver high performance on a reasonably priced Linux box. The index is used to find areas of probable homology, which are then loaded into memory for a detailed alignment. Protein BLAT works in a similar manner, except with 4-mers rather than 11-mers. The protein index takes a little more than 2 gigabytes.

BLAT was written by [Jian Kent](#). Like most of Jim's software interactive use on the web server is free to all. Sources and executables to run batch jobs on your own server are available free for academic, personal, and non-profit purposes. Non-exclusive commercial licenses are also available. Contact Jim for details.

NCBI Home > Genomic Biology > Human

Search [LocusLink] for [Go]

## The Human Genome

A guide to online information resources

**Web Resources**

- BLAST.** Compare your sequence to the genome or its gene products.
- Cytogenetics.** A cytogenetic resource of FISH-mapped, sequence-tagged clones.
- dbSNP.** Database of SNPs and other genetic variations.
- e-PCR.** Check your sequence for STSs and view in genomic context.
- GEO.** Gene Expression Omnibus, a public repository for expression data.
- HomoloGene.** Putative homologies among human, mouse, rat, and zebrafish.
- Homology Map.** Blocks of conserved synteny between mouse and human.
- LocusLink.** Focal point for genes and associated information.
- OMIM.** Guide to genes and inherited disorders maintained by JHU and collaborators.
- RefSeq.** Reference sequences of

**Building an information infrastructure**

A challenge facing researchers today is the ability to piece together and analyze the multitudes of data currently being generated through the Human Genome Project. NCBI's Web site serves as an integrated, one-stop, genomic information infrastructure for biomedical researchers from around the world so that they may use this data in their research efforts. [More...](#)

**Working Draft Analysis Published**

- NLM Press Release
- NHGRI Press Release
- Interactive Tour of the Genome
- NCBI Genome Analysis Pipeline
- Nature (2/15/01) Human Genome Issue
- Science (2/16/01) Human Genome Issue

**Browse**

Genes

**MapViewer tips and tricks**

When browsing the genome using the new MapViewer, click on Display Settings to choose from several types of maps and . Below are three views of the BRCA2 locus using different display options. Click the image to see the full MapViewer display.

PRO0297 - av sv 13 PRO0297 protein

BRCA2 - av sv 13 sv 2.3 Breast cancer 2, e

BRCA2 - av sv 13 sv 2.3 Breast cancer 2, e

BRCA2 - av sv 13 sv 2.3 Breast cancer 2, e

**Genes & Disease**

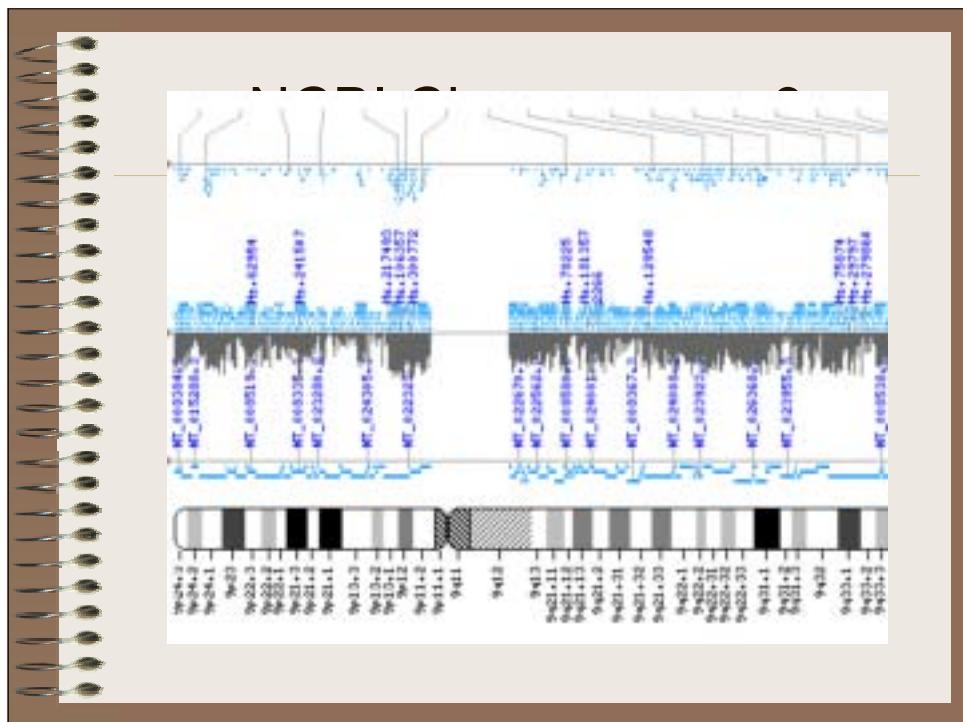
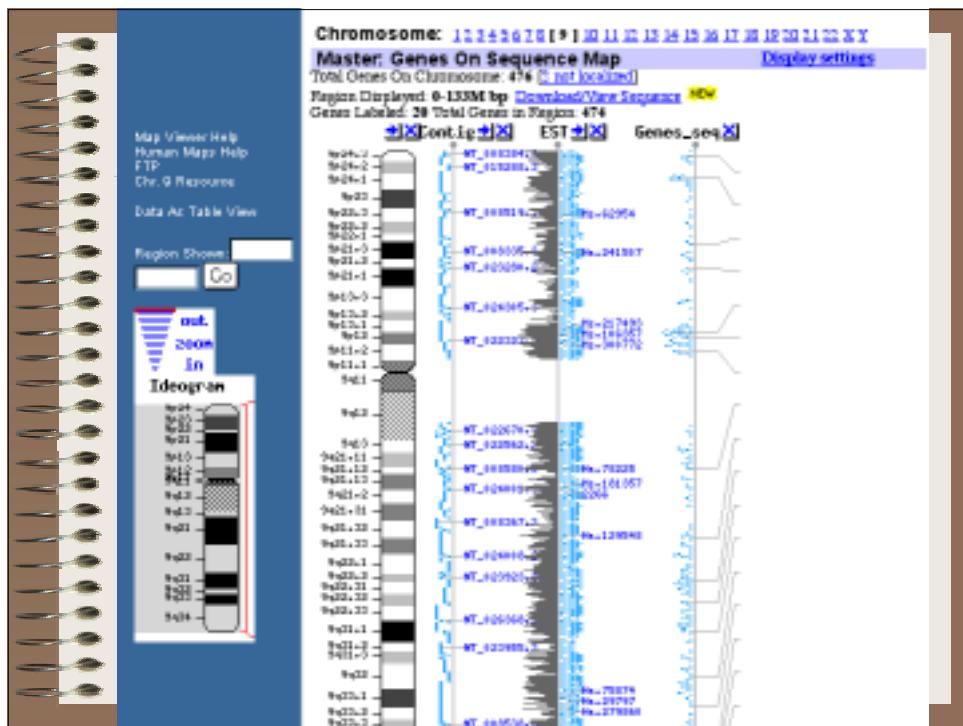
G&D. Selected gene stories for students and the public.

# NCBI Display options

A. Genes

B. Variations, juxtaposed with genes

C. Several STS maps, juxtaposed with genes



**About Ensembl v1.0**

Ensembl is a joint project between EMBL - EBI and The Sanger Centre to develop a software system which produces and maintains automatic annotation on eukaryotic genomes. Ensembl is primarily funded by the Wellcome Trust.

**With Ensembl you can ...**

- Download all data, free, without constraints
- Search the DNA from the human genome
- Browse chromosomes
- Find genes, SNPs and mouse genome matches
- Look for proteins and protein families

**How Do I ...**

- Find genomic sequences similar to my protein sequence?
- Look up a positional marker and examine candidate disease genes in the region?
- Find the expression profile of a gene?
- More...

**Ensembl provides ...**

- Identification of 90% of known human genes in the genome sequence
- Prediction of 10,000 additional genes, all with supporting evidence
- Connections to other resources worldwide, leveraging many public genomic databases and tools
- This website, www.ensembl.org, facilitates public access to this data by offering a web based genomic browser.

**Browse a Chromosome**

1 2 3 4 5 6 7 8 9 10 11 12  
13 14 15 16 17 18 19 20 21 22 X Y

**Ensembl Links**

- News
- Download
- BLAST
- SSAHA
- Docs
- Dev

New: [Ensembl Mouse server](#)

**Help**

Click on any help icon to pop up a context-sensitive help window.

**Marker**  **Lookup** [e.g. [RH9632, DIS2806](#)] **Help**

**Chromosome 9**

Known Ensembl Genes: 565 SNPs: 45992  
Novel Ensembl Genes: 478 Length: 141263275 bp

**Change Chromosome**

Chromosome:  **Lookup**

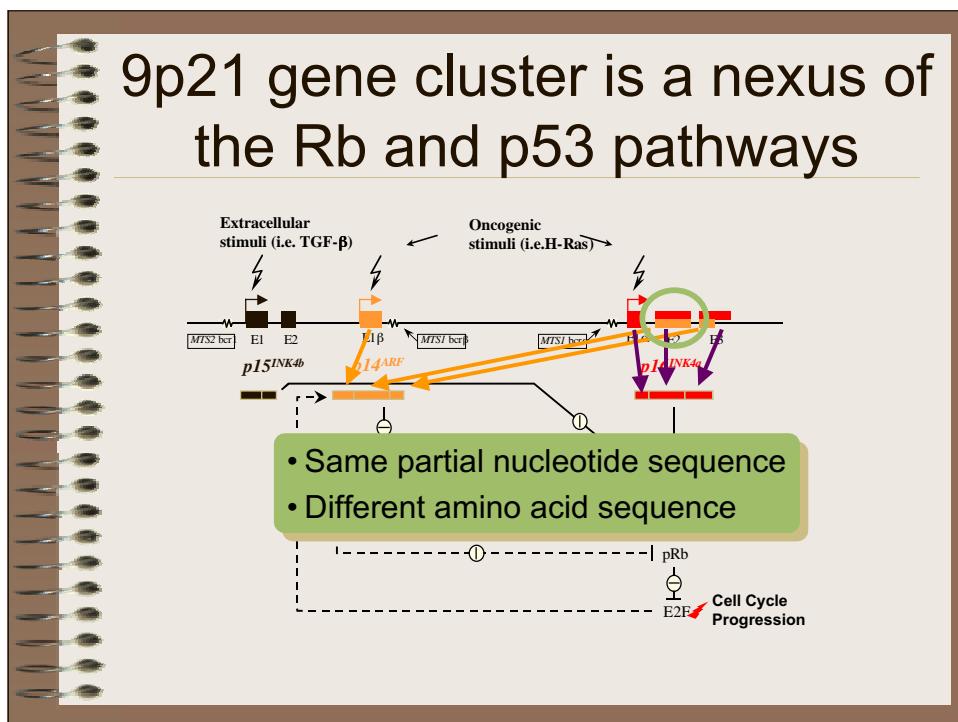
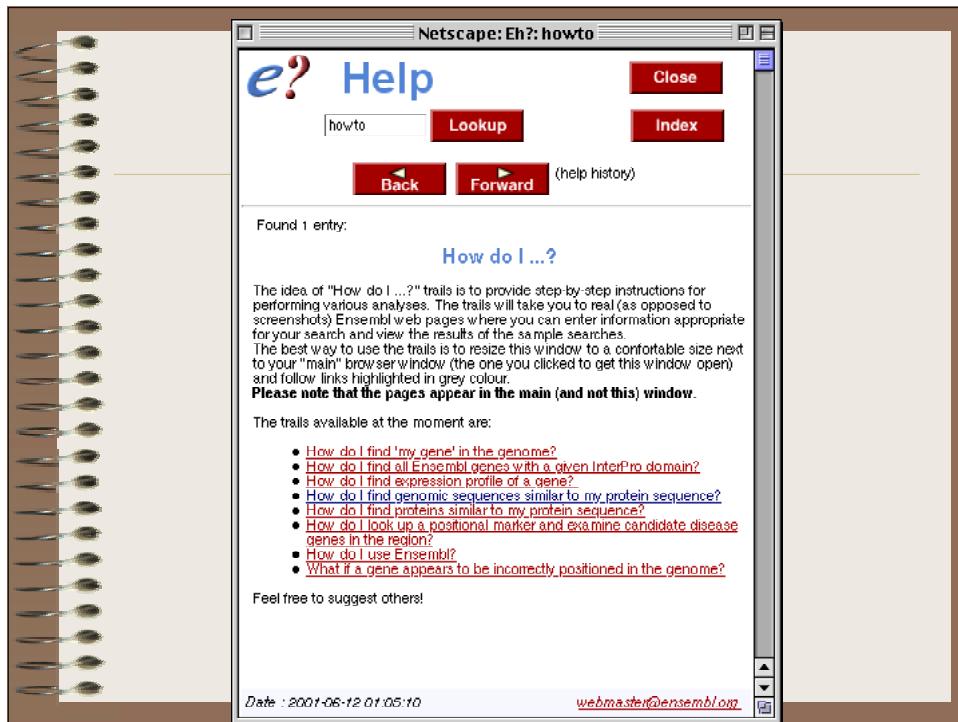
**Jump to Contigview**

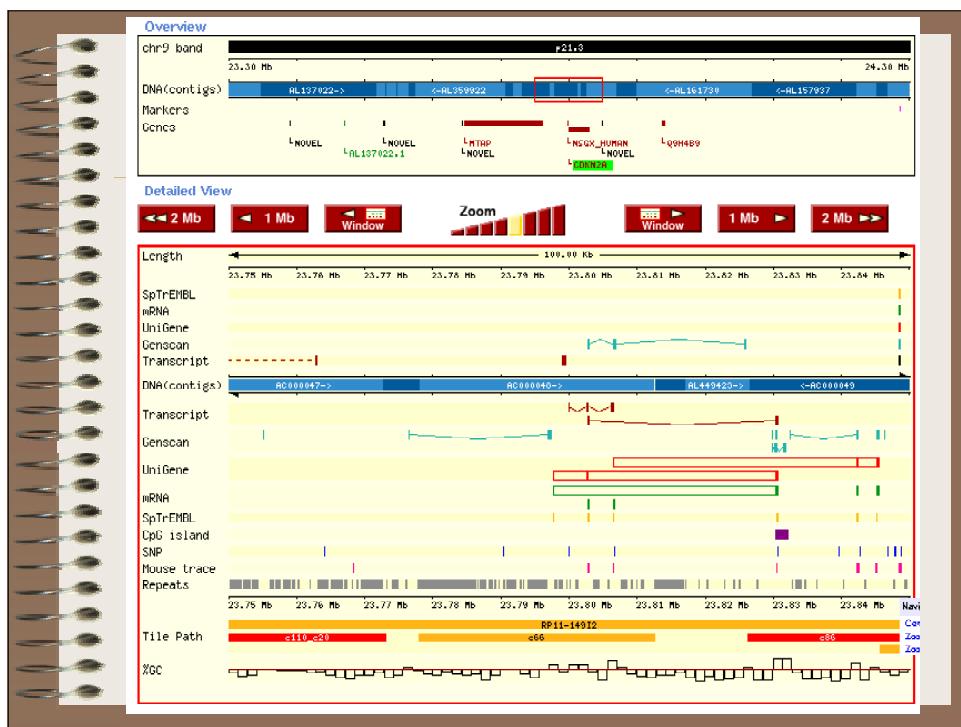
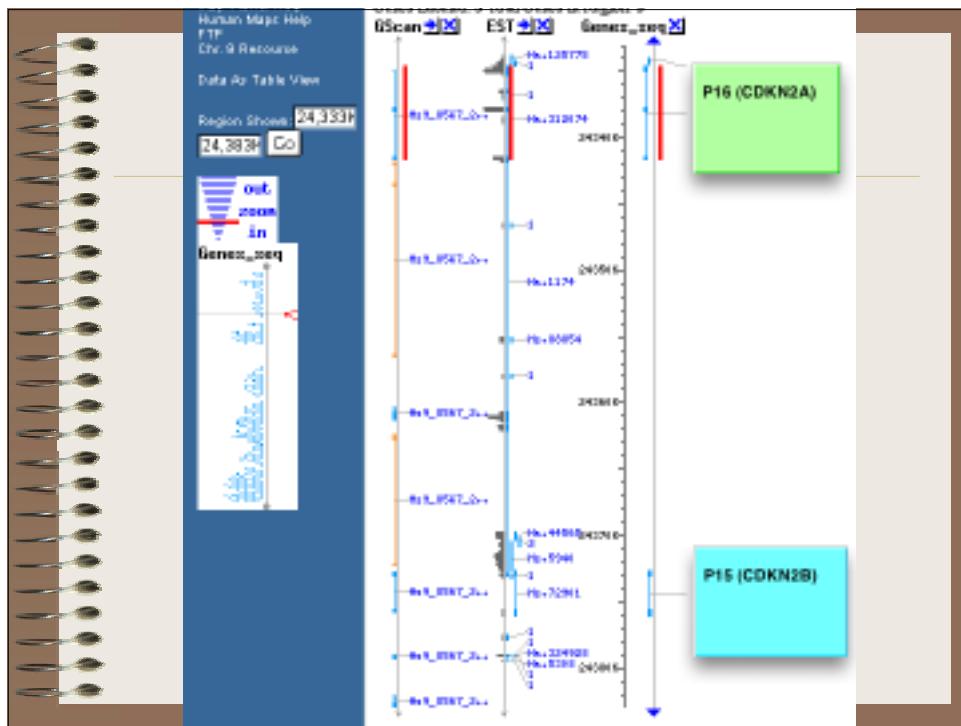
Click anywhere on the chromosome ideogram or one of the feature distribution plots to jump to a contig-level view of features at that point.  
Alternatively, you can jump to contigview between any two landmark markers on this chromosome.

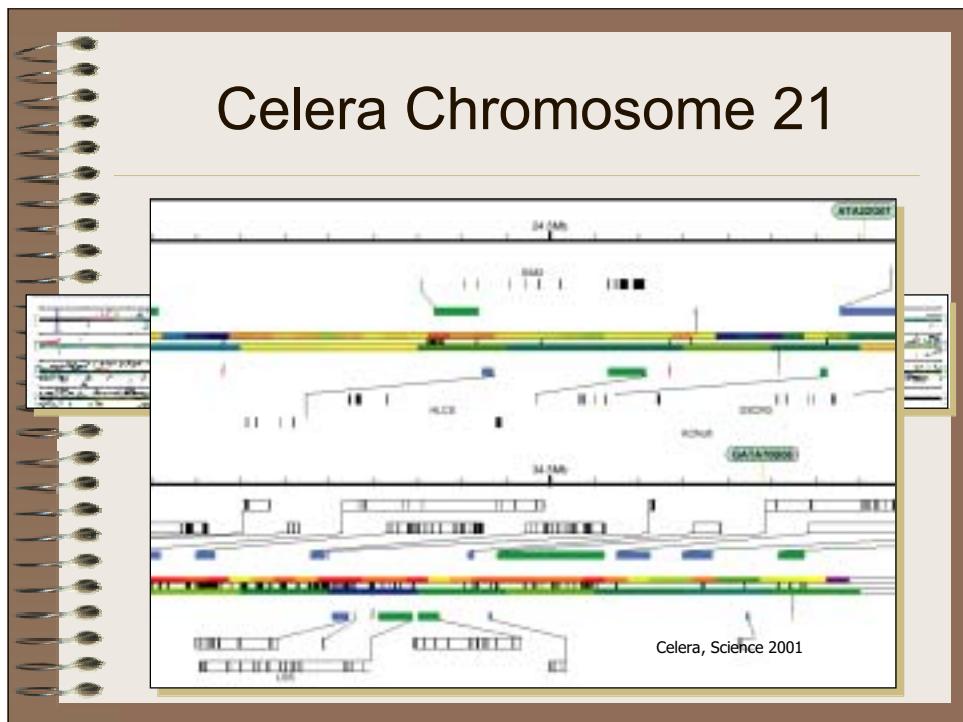
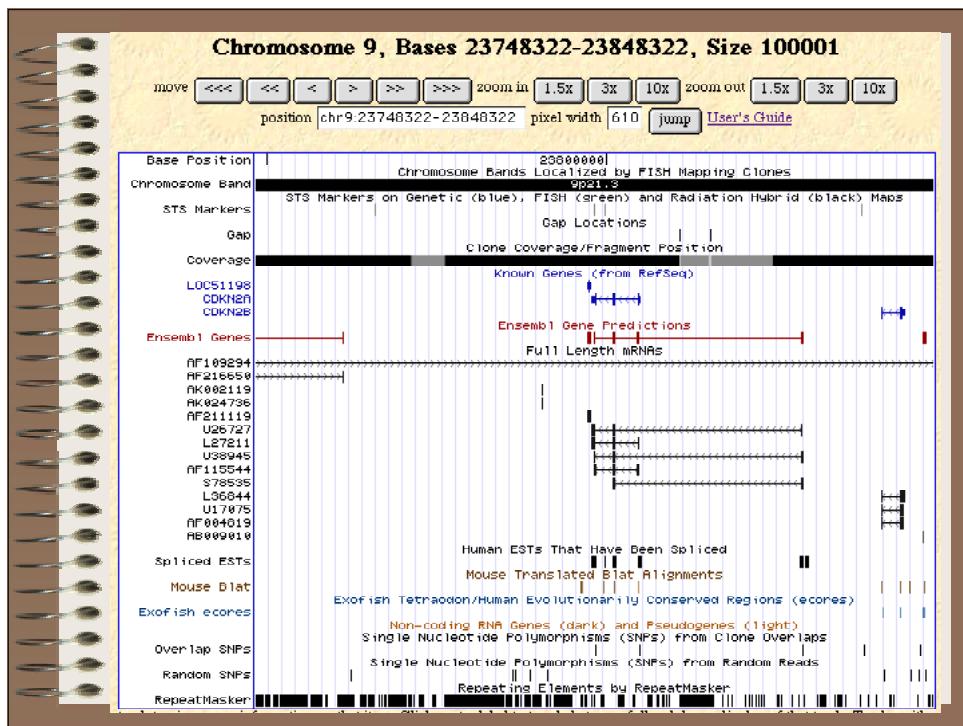
Between:  and  **Lookup**

**OMIM Diseases**

[Browse OMIM Diseases](#) on this chromosome.







# Beyond the Genome

This is the ExPASy (Expert Protein Analysis System) proteomics server of the [Swiss Institute of Bioinformatics](#) (SIB). This server is dedicated to the analysis of protein sequences and structures as well as 2-D PAGE ([Disclaimer](#))

[Announcements] [Job opening] [Mirror Sites]

Databases	Tools and Software Packages
<ul style="list-style-type: none"><li>• <a href="#">SWISS-PROT and TrEMBL</a> - Protein sequences</li><li>• <a href="#">PROSITE</a> - Protein families and domains</li><li>• <a href="#">SWISS-2DPAGE</a> - Two-dimensional polyacrylamide gel electrophoresis</li><li>• <a href="#">SWISS-3DIMAGE</a> - 3D images of proteins and other biological macromolecules</li><li>• <a href="#">SWISS-MODEL Repository</a> - Automatically generated protein models</li><li>• <a href="#">CD40Lbase</a> - CD40 ligand defects</li><li>• <a href="#">ENZYME</a> - Enzyme nomenclature</li><li>• <a href="#">SeqAnalRef</a> - Sequence analysis bibliography references</li><li>• <a href="#">Links to many other molecular biology databases</a></li></ul>	<ul style="list-style-type: none"><li>• <b>Proteomics tools</b><ul style="list-style-type: none"><li>◦ Identification and characterization</li><li>◦ DNA -&gt; Protein</li><li>◦ Similarity searches</li><li>◦ Pattern and profile searches</li><li>◦ Post-translational modification prediction</li><li>◦ Primary structure analysis</li><li>◦ Secondary structure prediction</li><li>◦ Tertiary structure</li><li>◦ Transmembrane regions detection</li><li>◦ Alignment</li></ul></li><li>• <a href="#">Molanal 3</a> - Software for 2-D PAGE analysis</li><li>• <a href="#">SWISS-MODEL</a> - Automated knowledge-based protein modeling server</li><li>• <a href="#">Swiss-PdbViewer</a> - Macintosh/PC tool for structure display and analysis</li><li>• <a href="#">Boehringer Mannheim's Biochemical Pathways</a></li></ul>
Education and services	Documentation
<ul style="list-style-type: none"><li>• <a href="#">The ExPASy FTP server</a></li><li>• <a href="#">Swiss-Shop</a> - automatically obtain (by email) new sequence entries relevant to your field(s) of interest</li><li>• <a href="#">Master Degree in Bioinformatics</a></li><li>• <a href="#">2-D PAGE training</a> - attend a one-week course in Geneva</li><li>• <a href="#">SWISS-2DSERVICE</a> - get your 2-D Gels performed according to Swiss standards</li></ul>	<ul style="list-style-type: none"><li>• <a href="#">What's New on ExPASy</a></li><li>• <a href="#">SWISS-FLASH electronic bulletins</a></li><li>• <a href="#">SWISS-PROT documents</a></li><li>• <a href="#">How to create HTML links to ExPASy</a></li><li>• <a href="#">Complete table of available documents</a></li></ul>
Links to lists of molecular biology resources	Links to some major molecular biology servers
<ul style="list-style-type: none"><li>• <a href="#">Amos' WWW links</a> - The ExPASy list of Biomolecular servers</li><li>• <a href="#">BioHunt</a> - Search the internet for molecular biology information</li><li>• <a href="#">WORLD-2DPAGE</a> - Links to 2-D PAGE database servers and 2-D PAGE related servers and services</li></ul>	<ul style="list-style-type: none"><li>• <a href="#">European Bioinformatics Institute (EBI)</a></li><li>• <a href="#">National Center for Biotechnology Information (NCBI)</a></li><li>• <a href="#">Japanese GenomeNet</a></li><li>• <a href="#">Australian National Genomic Information Service</a></li></ul>

# Physical Properties

## Prediction of Physical Properties

- Compute pi/MW <http://www.expasy.ch/tools/pi tool.html>
- MOWSE <http://srs.hgmp.mrc.ac.uk/cgi-bin/mowse>
- PeptideMass <http://www.expasy.ch/tools/peptide-mass.html>
- TGREASE <ftp://ftp.virginia.edu/pub/fa sta/>
- SAPS <http://www.isrec.isb-sib.ch/software/SAPSform.html>

## Prediction of Protein Identity Based on Composition

- AACompIdent <http://www.expasy.ch/tools/aacomp/>
- AACompSim <http://www.expasy.ch/tools/aacsim/>
- PROPSERCH <http://www.embl-heidelberg.de/prs.html>

## Motifs and Patterns

- BLOCKS <http://blocks.fhcrc.org>
- Pfam <http://www.sanger.ac.uk/Software/Pfam/>
- PRINTS <http://www.bioinf.man.ac.uk/dbbrowser/PRINTS/PRINTS.html>
- ProfileScan <http://www.isrec.isb-sib.ch/software/PFSCANform.html>

# Protein Structure

## Predictor of Secondary Structure and Folding Classes

- nnpredict [http://www.cmpharm.ucsf.edu/\\_nomi/nnpredict.html](http://www.cmpharm.ucsf.edu/_nomi/nnpredict.html)
- PredictProtein <http://www.embl-heidelberg.de/predictprotein/>
- SOPMA <http://pbil.ibcp.fr/>
- Jpred <http://jura.ebi.ac.uk:8888/>
- PSIPRED <http://insulin.brunel.ac.uk/psipred>
- PREDATOR <http://www.embl-heidelberg.de/predator/predatorinfo.html>

## Prediction of Specialized Structures or Features

- COILS <http://www.ch.embnet.org/software/COILSform.html>
- MacStripe <http://www.york.ac.uk/depts/biol/units/coils/mstr2.html>
- PHDtopology <http://www.embl-heidelberg.de/predictprotein>
- SignalP <http://www.cbs.dtu.dk/services/SignalP/>
- TMpred <http://www.isrec.isb-sib.ch/ftp-erver/tmpred/www/TMPREDform.html>

## Structure Prediction

- DALI <http://www2.ebi.ac.uk/dali/>
- Bryant-Lawrence <ftp://ncbi.nlm.nih.gov/pub/pkb/>
- FSSP <http://www2.ebi.ac.uk/dali/fssp/>
- UCLA-DOE <http://fold.doe-mbi.ucla.edu/Home>
- SWISS-MODEL <http://www.expasy.ch/swissmod/SWISS-MODEL.html>
- TOPITS <http://www.embl-heidelberg.de/predictprotein>